

Speech recognition method and device for carrying out the method

Patent Number: EP0817167
Publication date: 1998-01-07
Inventor(s): SCHEPPACH FRANK (DE)
Applicant(s): DAIMLER BENZ AEROSPACE AG (DE)
Requested Patent: ☐ EP0817167, A3
Application Number: EP19970110167 19970621
Priority Number(s): DE19961025294 19960625
IPC Classification: G10L3/00
EC Classification: G10L15/20
Equivalents: ☐ DE19625294

Abstract

The method involves passing the input speech through a preprocessor into a speech processing unit, where it is subjected to a speech recognition process. The preprocessor extracts characteristics from the input speech, segments and classifies it according to the energy content, and only passes the segments whose energy content exceeds a defined or adaptively determined threshold. Each segment whose energy content exceeds the threshold is finally investigated, and classified as to whether it is derived from a speech or non-speech signal. Only those segments derived from speech signals are processed further.

Data supplied from the esp@cenet database - I2

Best Available Copy

THIS PAGE BLANK (USPTO)

①9 BUNDESREPUBLIK

DEUTSCHLAND

⑫

Offenlegungsschrift

⑩

DE 196 25 294 A 1

⑤1

Int. Cl.⁶:

G 10 L 7/08



DEUTSCHES

PATENTAMT

⑳ Aktenzeichen: 196 25 294.6
 ㉔ Anmeldetag: 25. 6. 96
 ㉕ Offenlegungstag: 2. 1. 98

⑦1 Anmelder:

Daimler-Benz Aerospace Aktiengesellschaft, 81663
 München, DE

⑦2 Erfinder:

Scheppach, Frank, Dipl.-Ing., 89231 Neu-Ulm, DE

⑤4 Spracherkennungsverfahren und Anordnung zum Durchführen des Verfahrens

⑤7 Die Erfindung betrifft ein Spracherkennungsverfahren, bei dem eingehende Sprachsignale über eine Vorverarbeitungseinheit einer Spracherkennungseinheit zugeleitet werden und dort einem Spracherkennungsprozeß unterworfen werden, wobei mittels der Vorverarbeitungseinheit die eingehenden Sprachsignale einer Merkmalsextraktion, einer Segmentierung und einer Klassifizierung nach dem Energiegehalt unterzogen werden und wobei nur diejenigen Segmente weiterverarbeitet werden, deren Energiegehalt einen vorgegebenen oder adaptiv ermittelten Energie-Schwellenwert überschreitet.

Um eine möglichst robuste Spracherkennung durchführen zu können und um den eigentlichen Spracherkennung zu entlasten, wird nach der Erfindung vorgeschlagen, daß diejenigen Segmente, deren Energiegehalt den vorgegebenen oder adaptiv ermittelten Energie-Schwellenwert überschreiten, anschließend daraufhin untersucht und klassifiziert werden, ob sie aus einem Sprachsignal oder einem Nicht-Sprachsignal abgeleitet sind, und daß nur die als aus einem Sprachsignal abgeleitet klassifizierten Segmente weiter verarbeitet werden.

Dementsprechend wird für die Anordnung zum Durchführen des Verfahrens, die mit einer Merkmalsextraktionseinheit, einem dieser Einheit nachgeschalteten Energiedetektor, einer diesem Detektor nachgeschalteten Vektorquantisierungseinheit und einer dieser Einheit nachgeschalteten Klassifizierungseinheit ausgerüstet ist, vorgeschlagen, daß zwischen Energiedetektor und Vektorquantisierungseinheit

Die folgenden Angaben sind den vom Anmelder eingereichten Unterlagen entnommen

BUNDESDRUCKEREI 10. 97 702 081/392

11/22

Beschreibung

2

Die Erfindung betrifft ein Spracherkennungsverfahren gemäß Oberbegriff des Patentanspruchs 1 sowie eine Anordnung zum Durchführen des Verfahrens gemäß Oberbegriff des Patentanspruchs 7.

Es wurden bereits Spracherkennungsverfahren vorgeschlagen, bei denen eingehende Sprachsignale über eine Vorverarbeitungseinheit einer Spracherkennungseinheit zugeleitet werden und erst dort dem eigentlichen Spracherkennungsprozeß unterworfen werden. In der Vorverarbeitungseinheit werden die eingehenden Sprachsignale einer Merkmalsextraktion, einer Segmentierung und einer Klassifizierung nach dem Energiegehalt unterzogen. Bei diesem Verfahren werden nur diejenigen Segmente in der Spracherkennungseinheit weiterverarbeitet, deren Energiegehalt einen vorgebbaren oder adaptiv ermittelten Energie-Schwellenwert überschreitet. Mit dieser Maßnahme sollen störende Nebengeräusche wie z. B. Atemgeräusche, Papiergeknister, Tastaturklappern, Geräusche von Maschinen und Geräten, die während der Spracheingabe oder des Sprachdialogs auftreten, ausgeblendet werden. Akustische Signale, die den Energie-Schwellenwert überschreiten, werden in der Spracherkennungseinheit daraufhin untersucht, welchem Wort bzw. welchen Wortfolgen sie am wahrscheinlichsten entsprechen.

Dabei vergleicht die Spracherkennungseinheit die eingehenden Signale mit abgespeicherten Referenzsignalen (die bestimmten Worten entsprechen) und stellt dasjenige Referenzsignal fest, das die größte Übereinstimmung mit dem zu klassifizierenden Eingangssignal aufweist und das auch aufgrund einer Plausibilitätsprüfung von Bedeutungsinhalt des zugehörigen Wortes her nicht auszuschließen ist.

Da bei diesem Verfahren die Segmentierung und Weiterleitung an die Spracherkennungseinheit nur in Abhängigkeit von der Signalenergie erfolgt, werden neben den echten Sprachsignalen auch laute Nebengeräusche der Spracherkennungseinheit zugeführt und dort zwangsweise dem Spracherkennungsprozeß unterzogen mit der Folge, daß störenden lauten Nebengeräuschen, die den vorgegebenen Energie-Schwellenwert überschreiten in der Spracherkennungseinheit zwangsweise ein Wort bzw. eine Wortfolge zugeordnet wird.

Dies führt entweder dazu, daß die Spracherkennungseinheit aufgrund einer nachfolgenden Plausibilitätsprüfung zu keinem Ergebnis kommt und dann das Gesamt-signal (echtes Sprachsignal und störendes Nebengeräusch) als unsinnig verwirft, oder dazu, daß dem echten Sprachsignal im Erkennungsprozeß ein falscher Bedeutungsinhalt bzw. ein falsches Wort oder eine falsche Wortfolge zugeordnet, deren Wahl zu Fehlern in der weiteren Verarbeitung dieser Worte bzw. Wortfolgen führt.

Dies führt beispielweise dazu, daß bei einem per Spracheingabe zu bedienenden Telefon in einem Kraftfahrzeug (Kfz) die gewünschte Telefonverbindung nicht zustande kommt, da das Spracherkennungssystem des Telefons die eingesprochenen Ziffern der Telefonnummer des Teilnehmers, den der Telefonbediener anrufen möchte, wegen störender Nebengeräusche im Kfz falsch interpretiert und aufgrund der inkorrekt erkannten Ziffern entweder keine oder eine falsche Telefonverbindung herstellt.

Die Bedienung eines solchen Spracherkennungssystems erfordert eine hohe Disziplin des Sprechers bei der Spracheingabe sowie eine relativ ruhige Umgebung.

Um störende Nebengeräusche ausblenden zu können, wurde bereits vorgeschlagen, all die Nebengeräusche, die mit einer gewissen Wahrscheinlichkeit auftreten, als Referenzmuster für Nichtworte dem Spracherkennungssystem zur Verfügung zu stellen.

Mit einem solchen Verfahren können diese Nebengeräusche zwar im Spracherkennungsprozeß als solche erkannt und eliminiert werden. Neu auftretende und nicht als Referenzmuster gespeicherte Nebengeräusche durchlaufen jedoch auch den gesamten Erkennungsprozeß und werden je nachdem, welches der gespeicherten Referenzmuster für Sprache oder Nichtworte ihnen am ähnlichsten ist, entweder (im günstigsten Fall) einem der abgespeicherten Nichtworte zugeordnet mit der Folge, daß sie nicht weiter berücksichtigt werden, oder (im ungünstigsten Fall) einem der abgespeicherten Worte bzw. Wortfolgen zugeordnet mit der Folge, daß sie weiterhin im Spracherkennungsprozeß berücksichtigt werden, was zu unsinnigen Ergebnissen führen kann.

Die Aufgabe der Erfindung besteht darin, zum einen ein Spracherkennungsverfahren zu schaffen, das eine möglichst hohe Robustheit aufweist und möglichst unempfindlich ist gegenüber störenden Nebengeräuschen bei der Spracheingabe, sowie zum anderen eine Anordnung zum Durchführen des Verfahrens zu schaffen, die möglichst einfach im Aufbau ist.

Die erfindungsgemäße Lösung der Aufgabe ist im Hinblick auf das zu schaffende Verfahren durch die kennzeichnenden Merkmale des Patentanspruchs 1 wiedergegeben und im Hinblick auf die zu schaffende Anordnung durch die kennzeichnenden Merkmale des Patentanspruchs 7.

Ein wesentlicher Vorteil der erfindungsgemäßen Lösung ist darin zu sehen, daß störende Nebengeräusche bereits in der Vorverarbeitungseinheit zuverlässig als solche erkannt werden und gar nicht erst der Spracherkennungseinheit zugeleitet werden. Dadurch wird die Folge, daß der Rechneraufwand bzw. die Rechenzeit erheblich verringert werden können.

Zudem entfällt die Notwendigkeit, typisch auftretende Nebengeräusche als Nichtwort-Referenzmuster in der Spracherkennungseinheit abspeichern zu müssen.

Da sowohl typische als auch atypische Nebengeräusche in der Vorverarbeitungseinheit als solche erkannt und eliminiert werden, arbeitet das erfindungsgemäße Verfahren mit einer sehr hohen Erkennungswahrscheinlichkeit für die eingesprochenen Sprachsignale. Die nach dem Verfahren arbeitende erfindungsgemäße Anordnung zeichnet sich daher durch ihre hohe Robustheit und Erkennungszuverlässigkeit selbst in erheblich nebengeräuschbehafteten Sprachumgebungen wie z. B. in Kraftfahrzeugen oder in Produktionsstätten bzw. Maschinenhallen aus.

Im folgenden wird die Erfindung anhand der Figuren näher erläutert. Es zeigen

Fig. 1 das Blockschaltbild eines Teils eines bevorzugten Ausführungsbeispiels der erfindungsgemäßen Anordnung zum Durchführen des erfindungsgemäßen Verfahrens.

Fig. 2 den zeitlichen Verlauf eines Sprachsignals mit eingezeichneten Segmentierungsintervallen für die der Spracherkennungseinheit zugeleiteten Sprachsignalteile bei einem Verfahren nach dem Stand der Technik (Fig. 2a) und bei einer bevorzugten Ausführungsform des erfindungsgemäßen Verfahrens gemäß Fig. 3 (Fig. 2b).

Fig. 3 das Flußdiagramm einer bevorzugten Ausführungsform.

rungsform des erfindungsgemäßen Verfahrens.

Das in Fig. 1 gezeigte Spracherkennungssystem besteht aus einer Spracheingabeeinheit (z. B. einem Mikrophon) 1, das ausgangsseitig an eine Merkmalsextraktionseinheit 2 angeschlossen ist. Die Merkmalsextraktionseinheit 2 ist ausgangsseitig an einen Energiedetektor 3 angeschlossen, dem ein Nichtwortdetektor 6 nachgeschaltet ist. Diesem wiederum ist eine Vektorquantisierungseinheit 4 nachgeschaltet, an die sich eine Klassifizierungseinheit 5 anschließt. Die sich daran anschließende — an sich bekannte — Signalverarbeitungs- und -ausgabereinheit des Spracherkennungssystems ist nicht gezeigt und bildet auch nicht den Gegenstand der vorliegenden Erfindung. Während die Merkmalsextraktionseinheit 2, der Energiedetektor 3 und der Nichtwortdetektor 6 zur Vorverarbeitungseinheit des Spracherkennungssystems gehören, sind die Vektorquantisierungseinheit 4 und die Klassifizierungseinheit 5 bereits Teil des eigentlichen Spracherkenners.

Die Einstellung und Funktionsweise der einzelnen in Fig. 1 gezeigten Funktionseinheiten (mit Ausnahme des Nichtwortdetektors 6) sind an sich bekannt und werden daher hier an dieser Stelle nicht weiter erläutert.

Der Nichtwortdetektor 6 wird off-line über bestimmte vorgegebene Parameter 7 eingestellt. Zum einen wird eine Initialisierung für Sprachsegmente durchgeführt (Abgleich mit verschiedenen Sprechern bzw. Detektion auf Vokale usw.), zum anderen erfolgt in dieser Phase anschließend eine Feinjustierung mit Nichtsprachsignalen (Geräuschen), bei der die Ausblendung solcher Signale eingestellt wird.

Zur Funktionsweise der Anordnung:

Über die Spracheingabeeinheit 1 eingehende Signale (Sprachsignale und Geräusche) werden in der Merkmalsextraktionseinheit 2 in Merkmalsvektoren zerlegt, die anschließend im Energiedetektor 3 daraufhin untersucht werden, ob sie einen vorgegebenen oder adaptiv ermittelten Energie-Schwellenwert überschreiten oder nicht.

Nur diejenigen Segmente, deren Energiegehalt diesen ersten Energie-Schwellenwert überschreiten, werden dem nachfolgenden Nichtwortdetektor zugeleitet.

Dort werden sie daraufhin untersucht und klassifiziert, ob sie aus einem Sprachsignal oder aus einem Nichtsprachsignal abgeleitet sind. Während die als von einem Nicht-Sprachsignal abgeleitet klassifizierten Segmente bei der weiteren Signalverarbeitung ausgeblendet, d. h. eliminiert werden, werden die als von einem Sprachsignal abgeleiteten klassifizierten Segmente der Vektorquantisierungseinheit 4, d. h. der eigentlichen Spracherkennungseinheit zugeleitet und dort weiterverarbeitet.

Wie drastisch sich diese Ausblendung der Nichtsprachsignale auf die weitere Verarbeitung auswirken kann, ist beispielhaft in Fig. 2 gezeigt.

Dort ist der zeitliche Verlauf der Amplitude eines Sprachsignals gezeigt, bei dem die Ziffernfolge "null, eins, zwei, drei, vier, fünf" als Worte gesprochen eingegeben worden sind (in der Figur sind unter die entsprechenden (echten) Sprachsignale die Ziffern ausgeschrieben worden).

Die Ziffern wurden von Atemgeräuschen unterbrochen in das System eingegeben, die in Fig. 2 im zeitlichen Verlauf des Sprachsignals ebenfalls sichtbar sind.

In Fig. 2a ist ferner die Energiesegmentierung gezeigt, wie sie in einem herkömmlichen Spracherkennung ohne Nicht-Sprachsignal-detektor ("Nichtwortdetektor") erfolgt. Alle Signalanteile, die den vorgegebenen

oder adaptiv ermittelten Energie-Schwellenwert überschreiten, sind in in der Figur in Rechtecke eingefasst worden. Hierzu zählen neben den tatsächlich eingesprochenen Ziffern 0 bis 5 auch die Atemgeräusche zwischen den Ziffereingaben. Dem Spracherkennung standen in diesem Beispiel als Referenz-Vokabularium nur die Ziffern 0 bis 9 zur Verfügung. Das heißt alle dem Spracherkennung zugeleiteten Signale (egal ob Sprachsignal oder Atemgeräusch), die den vorgegebenen oder adaptiv ermittelten Energie-Schwellenwert überschritten hatten, wurden im eigentlichen Spracherkennung zwangsweise diesen als Referenzmuster abgespeicherten Ziffern zugeordnet.

Dementsprechend wurden zwar die gesprochenen Ziffern 0 bis 5 richtig erkannt. Jedoch wurden die Atemgeräusche nicht als solche erkannt, sondern es wurde ihnen jeweils eine Ziffer aus dem Vokabularium des Spracherkenners zugeordnet, die in der Figur jeweils als Zahl dargestellt worden ist. Bindestriche zwischen den Zahlen stehen hier für im Verbundwortmodus erkannte Signale.

In Fig. 2b ist der zeitliche Verlauf desselben Sprachsignals gezeigt. Hier wurde jedoch zusätzlich in der Vorverarbeitungseinheit ein Nicht-Sprachsignal-detektor eingesetzt, und zwar in einer Anordnung gemäß Fig. 1. Während bei diesem System die eingesprochenen Ziffern 0 bis 5 eindeutig als Sprachsignal erkannt und klassifiziert worden sind (in der Fig. 2a erkennbar an den Rechtecken über den entsprechenden Signalteilen), sind die Atemgeräusche eindeutig als Nicht-Sprachsignale erkannt und für die weitere Verarbeitung ausgeblendet worden. Als Sprachsignale erkannt wurden mit diesem System folgerichtig nur die eingesprochenen Ziffern 0 bis 5 (in Fig. 2b unter den Signalanteilen als Ziffer dargestellt).

Dieser Vergleich unterstreicht die Robustheit des erfindungsgemäßen Verfahrens gegenüber störenden Nebengeräuschen.

Die Funktionsweise des Nicht-Sprachsignal-detektors (4 in Fig. 1) wird im folgenden anhand des Flußdiagramms der Fig. 3 näher erläutert.

Das Zeitsignal des eingesprochenen Sprachsignals wird zunächst einem Signalrahmenbildungsprozeß unterworfen und anschließend in den Spektralbereich transformiert. Dies geschieht hier beispielhaft mit Hilfe einer Fast-Fourier-Transformation (FFT). Die transformierten Signale werden anschließend dem Energiedetektor zugeleitet. Dort wird geprüft, ob die Signale in ihrem Energiegehalt einen vorgegebenen oder — alternativ hierzu — einen adaptiv ermittelten Energie-Schwellenwert überschreiten oder nicht. Wird der Schwellenwert nicht überschritten, wird das untersuchte Signal im weiteren Verlauf der Signalverarbeitung nicht mehr berücksichtigt ("ausgeblendet"), und das nächstfolgende Signal wird dem Energiedetektor zugeleitet (diese Möglichkeit ist in der Fig. 3 durch den "nein"-Pfeil graphisch dargestellt). Überschreitet das Signal jedoch mit seinem Energiegehalt den ersten Energie-Schwellenwert, wird das Signal im Nicht-Sprachsignal-detektor wie folgt weiterverarbeitet.

Zunächst wird das Signalspektrum gemäß der an sich bekannten mel-Skala in Teilbereiche aufgeteilt. Unter dem Begriff mel-Skala versteht man eine nichtlineare Aufteilung des hörbaren Frequenzbereichs, die an die Hörcharakteristik des menschlichen Ohres angepaßt ist (vgl. hierzu z. B. Rabiner, L.; Bing-Hwang, J.: "Fundamentals of Speech Recognition" (Prentice Hall, Englewood Cliffs, New Jersey, 1993, Seiten 183 bis 186). Das

5 melskalierte Spektrum wird anschließend in Teilbändern mit $i = 0, 1, 2, \dots, N-1$ zusammengefaßt.

Danach werden in einem rekursiven Verfahren (dargestellt in Fig. 3 durch den zugehörigen "nein"-Pfeil) alle Teilbänder separat untersucht und der in ihnen jeweils vorhandene partielle Energiegehalt, bestehend aus dem Maximum der Energie, $E_{i, \max}$ des Teilbandes, normiert auf den Gesamtenergiegehalt E_i des Teilbandes i , festgestellt.

Wenn die normierten partiellen Energiegehalte aller Teilbänder i feststehen, wird im nächsten Verfahrensschritt die Anzahl der Teilbänder i ermittelt, deren partielle Energiegehalte einen vorgegebenen Parameter 1 überschreiten (in der Figur steht das Symbol # für den Begriff "Anzahl").

Überschreitet der partielle Energiegehalt von mehr als einem Teilband und von weniger als N Teilbändern den vorgegebenen Parameter 1, werden die Bedingungen für ein artikulatorisches Geräusch als erfüllt angesehen und das Signal weiterverarbeitet.

Wird diese Bedingung jedoch nur von einem Teilband i erfüllt (in Fig. 3 dargestellt durch den zugehörigen "nein"-Pfeil), wird auf der Basis aller Teilbänder $i = 0, \dots, N-1$ zusätzlich der durchschnittliche partielle Energiegehalt gebildet und mit einem vorgegebenen Parameter 2 verglichen.

Liegt dieser Durchschnittswert über dem Parameter 2 und der gleichzeitig individuelle partielle Energiegehalt von (im Beispiel) einem Teilband über dem Parameter 1, so werden die Bedingungen für ein artikulatorisches Geräusch als erfüllt angesehen und das Signal weiterverarbeitet.

Wird diese zweifache Bedingung (individueller partieller Energiegehalt von einem Teilband größer als Parameter 1 und gleichzeitig durchschnittlicher Energiegehalt aller Teilbänder größer als Parameter 2) nicht erfüllt, wird ebenso wie für den Fall, daß der individuelle partielle Energiegehalt von allen N Teilbändern i ($i = 0, 1, \dots, N-1$) jeweils größer als Parameter 1 ist, das untersuchte Signal als ein nichtartikulatorisches Signal angesehen und im weiteren Verlauf der Signalverarbeitung nicht mehr berücksichtigt ("ausgeblendet"). Anschließend wird der nächstfolgende Signalrahmen dem Energiedetektor zugeleitet (in Fig. 3 ist dieser Fall durch den zugehörigen "nein"-Pfeil graphisch dargestellt).

Liegt der individuelle partielle Energiegehalt von mindestens zwei Teilbändern i über dem Parameter 1 oder — bei Nichterfüllung dieser Bedingung — liegt zumindest der durchschnittliche partielle Energiegehalt von allen N Teilbändern i ($i = 0, 1, \dots, N-1$) über dem Parameter 2 sowie der partielle Energiegehalt von einem Teilband über dem Parameter 1, wird das Signal zusammen mit den nachfolgenden Signalen, die ebenfalls die zuvor beschriebene Energiegehaltsprüfung durchlaufen haben, daraufhin überprüft, ob die Anzahl der diese Bedingungen erfüllenden Signalrahmen eine vorgegebene Anzahl von Signalrahmen entsprechend einer vorgegebenen Zeitdauer überschreitet oder nicht.

Überschreitet diese Folge von Signalrahmen die vorgegebene Zeitdauer, werden die Bedingungen für ein Sprachsignal als erfüllt angesehen und die Folge von Signalrahmen der weiteren Verarbeitung (Klassifikationsprozeß) zugeleitet.

Überschreitet diese Folge von Signalen die vorgegebene Zeitdauer jedoch nicht, werden die Bedingungen für ein Sprachsignal nicht als erfüllt angesehen und nach dem allgemein bekannten First-In-First-Out-Prinzip das

6 älteste dieser Signale von der weiteren Signalverarbeitung ausgeschlossen sowie anschließend der nächste Signalrahmen dem Energiedetektor zugeleitet (in Fig. 3 durch den "nein"-Pfeil graphisch dargestellt).

Das Verfahren zeichnet sich durch folgende Vorteile aus:

- Aufwendiges Training zur Modellierung von stimmlosen, artikulatorischen und nichtartikulatorischen Störungen kann entfallen.
- Hierdurch nicht erforderliche Referenzmuster oder Modelle ersparen Speicherplatz.
- Einsparung von Rechenaufwand bei Nachverarbeitung zugunsten einer beschleunigten Endauswertung.
- Einstellung der Parameter erfordert nur wenige Sprecher für sprecherunabhängige Nichtwortdetektion.

20 Es versteht sich, daß die Erfindung nicht auf die dargestellten Ausführungsbeispiele beschränkt ist, sondern vielmehr auf weitere übertragbar ist.

Es lassen sich beispielsweise die Parameter für den Nichtwortdetektor adaptiv oder über eine vorgegebene Datenbasis individuell auf die jeweilige Anwendung bezogen einstellen.

25 Ferner ist es z. B. möglich, anstelle der FFT andere Transformationsverfahren einzusetzen, um das anstehende Zeitsignal in den Spektralbereich zu transformieren.

Auch ist es möglich, bei der Entscheidung, ob die Bedingungen für ein artikulatorisches Geräusch vorliegen, auf der Basis der durchschnittlichen und individuellen partiellen Energiegehalte die Mindestzahl j der Teilbänder, deren partieller Energiegehalt über dem Parameter 1 liegen muß, damit die Bedingungen für ein artikulatorisches Geräusch als erfüllt angesehen werden, individuell für den Anwendungszweck anzupassen (von z. B. $j = 2$ auf $j = 3$ oder 5 oder 10 usw.; allgemein $1 < j < N$). Entsprechend kann auch für den Fall, daß als zweite Bedingung der durchschnittliche partielle Energiegehalt bestimmt und mit einem zweiten Parameter verglichen werden muß, die Zahl k der Teilbänder, deren partieller Energiegehalt den ersten Parameter zumindest überschritten haben muß, individuell für den Anwendungszweck angepaßt werden ($1 \leq k < j$).

Ferner kann die Anzahl N der Teilbänder i an den individuellen Anwendungszweck angepaßt werden.

Schließlich ist es möglich, den Energie-Schwellenwert selbst an die jeweilige Anwendung anzupassen, d. h. zu erhöhen oder zu erniedrigen. Auch die vorgegebene Zeitdauer, die die einzelnen Signalfolgen erfüllen müssen, damit die Bedingungen für ein Sprachsignal als erfüllt angesehen werden, kann dementsprechend an den individuellen Anwendungszweck angepaßt, d. h. verlängert oder verkürzt werden.

Auch können anstelle der Energiedetektoren auch Pitchdetektoren eingesetzt werden (oder beide Detektorarten in Kombination betrieben werden).

Patentansprüche

1. Spracherkennungsverfahren, bei dem eingehende Sprachsignale über eine Vorverarbeitungseinheit einer Spracherkennungseinheit zugeleitet werden und dort einem Spracherkennungsprozeß unterworfen werden, wobei mittels der Vorverarbeitungseinheit die eingehenden Sprachsignale einer

Merkmalsextraktion, einer Segmentierung und einer Klassifizierung nach dem Energiegehalt unterzogen werden und wobei nur diejenigen Segmente weiterverarbeitet werden, deren Energiegehalt einen vorgegebenen oder adaptiv ermittelten ersten Energie-Schwellenwert überschreitet, dadurch gekennzeichnet, daß diejenigen Segmente, deren Energiegehalt den vorgegebenen oder adaptiv ermittelten Energie-Schwellenwert überschreiten, anschließend daraufhin untersucht und klassifiziert werden, ob sie aus einem Sprachsignal oder einem Nicht-Sprachsignal abgeleitet sind, und daß nur die als aus einem Sprachsignal abgeleitet klassifizierten Segmente weiterverarbeitet werden.

2. Spracherkennungsverfahren nach Anspruch 1, dadurch gekennzeichnet, — daß die Merkmalsextraktion eine Spektraltransformation beinhaltet und die spektral transformierten und in ihrem Energiegehalt den vorgegebenen oder adaptiv ermittelten Energie-Schwellenwert überschreitenden Segmente nach einer vorgegebenen mel-Skala in spektrale Bereiche aufgeteilt und anschließend die spektralen Bereiche in Teilbänder i mit $i = 0,1,2 \dots N-1$ zusammengefaßt werden;

— daß in allen Teilbändern i jeweils der als das auf den Gesamtenergiegehalt E_i des jeweiligen Teilbandes i normierte Maximum der Energie $E_{i,max}$ des jeweiligen Teilbandes i definierte partielle Energiegehalt des jeweiligen Teilbandes bestimmt wird;

— daß anschließend die partiellen Energiegehalte aller Teilbänder i mit einem vorgegebenen ersten Parameter (Parameter 1) verglichen werden und

— daß bei Überschreiten des ersten Parameters (Parameter 1) durch die partiellen Energiegehalte von j Teilbändern, $1 < j < N$, vorgegebene Bedingungen für das Vorliegen eines artikulatorischen Geräusches als erfüllt angesehen werden und das zugehörige Segment des Sprachsignals weiterverarbeitet wird.

3. Spracherkennungsverfahren nach Anspruch 2, dadurch gekennzeichnet,

— daß bei Segmenten, bei denen der partielle Energiegehalt nur von $j - k$ Teilbändern, $1 < j < N$ und $1 \leq k < j$, den vorgegebenen ersten Parameter (Parameter 1) überschreitet, anschließend der durchschnittliche partielle Energiegehalt aus allen partiellen Energiegehalten des jeweiligen Segments oder aus einem Teil dieser partiellen Energiegehalte des jeweiligen Segments gebildet und mit einem vorgegebenen zweiten Parameter (Parameter 2) verglichen wird und

— daß bei Überschreiten des zweiten Parameters (Parameter 2) durch den durchschnittlichen partiellen Energiegehalt die vorgegebenen Bedingungen für das Vorliegen eines artikulatorischen Geräusches als erfüllt angesehen werden und das zugehörige Segment des Sprachsignals weiterverarbeitet wird.

4. Spracherkennungsverfahren nach einem der Ansprüche 2 oder 3, dadurch gekennzeichnet, daß für den Fall, daß die partiellen Energiegehalte aller Teilbänder i den vorgegebenen ersten Parameter (1) überschreiten oder daß der partielle Energiegehalt von keinem der Teilbänder i den vorgegebenen

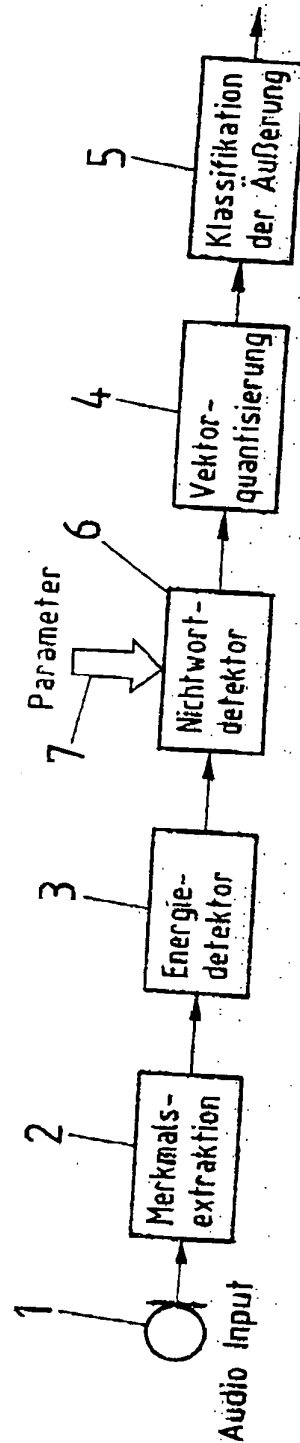
ersten Parameter (Parameter 1) überschreitet, oder für den Fall, daß zwar der partielle Energiegehalt von $j-k$ Teilbändern den ersten Parameter (Parameter 1) überschreitet, aber der durchschnittliche partielle Energiegehalt nicht den zweiten Parameter (Parameter 2) überschreitet, die vorgegebenen Bedingungen für das Vorliegen eines artikulatorischen Geräusches nicht als erfüllt angesehen werden und das zugehörige Segment von der weiteren Verarbeitung ausgeschlossen wird.

5. Spracherkennungsverfahren nach einem der vorhergehenden Ansprüche 2 bis 4, dadurch gekennzeichnet, daß $j = 2$ gewählt ist.

6. Spracherkennungsverfahren nach einem der Ansprüche 2 bis 5, dadurch gekennzeichnet, daß, sofern für eine vorgegebene Anzahl von aufeinander folgenden Segmenten die Bedingungen für das Vorliegen eines artikulatorischen Geräusches als erfüllt angesehen werden, diese Gruppe von Segmenten als Sprachsignal klassifiziert und weiterverarbeitet werden.

7. Anordnung zum Durchführen des Verfahrens nach einem der vorhergehenden Ansprüche, mit einer Merkmalsextraktionseinheit, einem dieser Einheit nachgeschalteten Energiedetektor, einer diesem Detektor nachgeschalteten Vektorquantisierungseinheit und einer dieser Einheit nachgeschalteten Klassifizierungseinheit, dadurch gekennzeichnet, daß zwischen Energiedetektor (3) und Vektorquantisierungseinheit (4) ein Nichtwortdetektor (6) geschaltet ist.

Hierzu 3 Seite(n) Zeichnungen



Parametereinstellung (offline):

- Initialisierung für Sprachsegmente (Abgleich mit verschiedenen Sprechern, Detektion auf Vokale)
- Feinjustierung mit Nichtsprache bis zum Ausblende effekt

FIG. 1

Nummer:
 Int. Cl. 6:
 Offenlegungstag:

DE 196 25 294 A1
 G 10 L 7/08
 2. Januar 1998

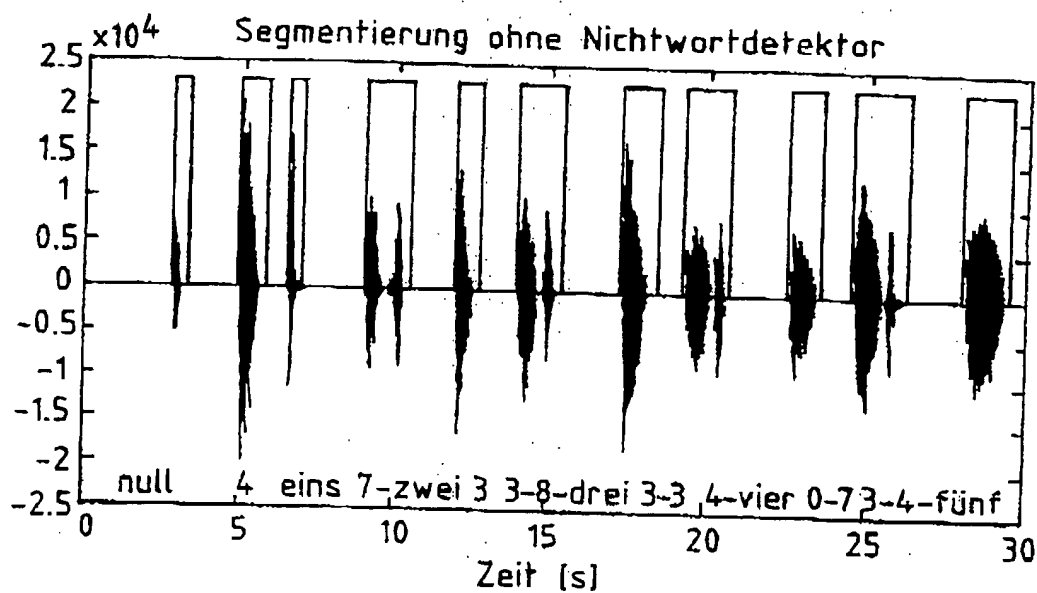


FIG. 2A

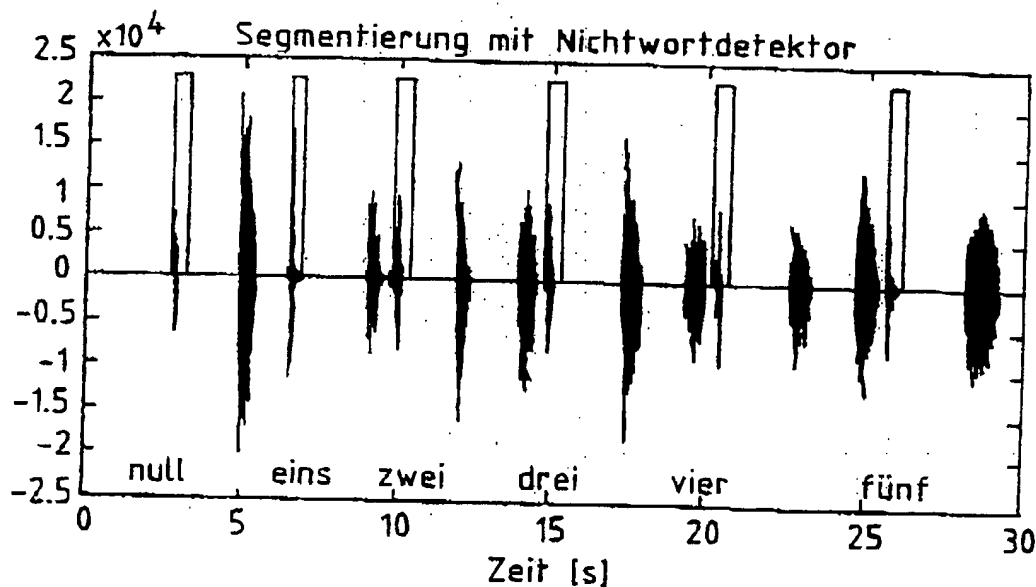


FIG. 2B

Nummer:
Int. Cl. 6:
Offenlegungstag:

DE 196 25 294 A1
G 10 L 7/08
2. Januar 1998

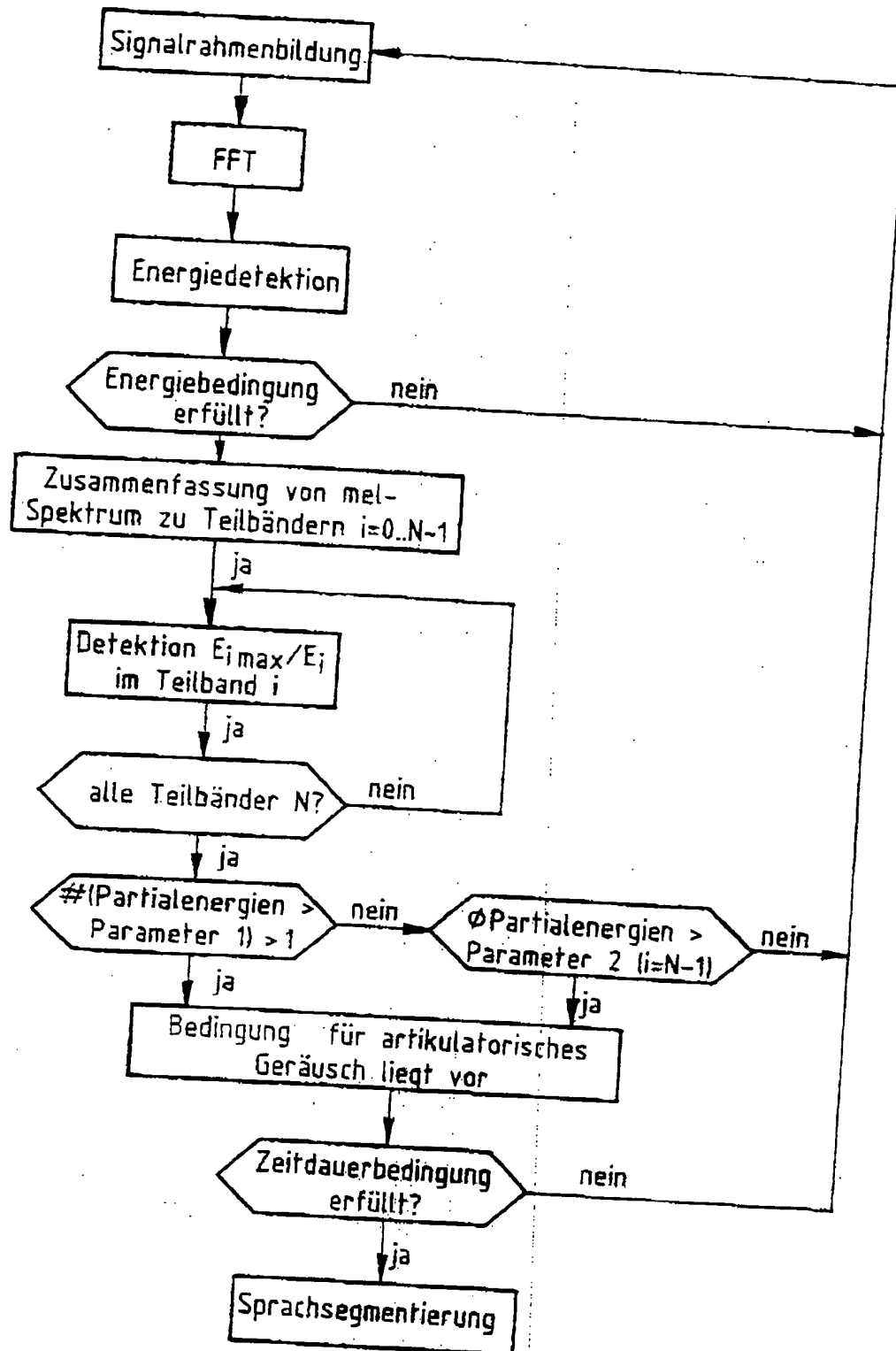


FIG. 3

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☒ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.

This Page Blank (uspto)